

STORM

流式数据处理

by xiaofei

概念

Storm is a free and open source distributed realtime computation system.

描述

Storm makes it easy to reliably process unbounded streams of data, doing for **realtime processing** what Hadoop did for **batch processing**. Storm is simple, can be used with any programming language, and is a lot of fun to use!

功能特点

便于集成

Spout已实现与Kestrel、RabbitMQ/AMQP、Kafka、JMS等的集成插件

开发简单易于使用

只需要实现Spout、 Bolt和topology,将topology提交到集群即可

水平扩展

可通过增加机器扩展Spout和Bolt的处理能力

保障数据被处理

- 记录级容错
- Transaction

容错

- **worker死掉了**:Supervisor将重启它，持续fail且无心跳到Nimbus，Nimbus将重新分配它到其它机器。
- **node死掉了**:task分配到的机器无响应或time out，Nimbus将分配到其它机器
- **Nimbus或Supervisor服务死掉了**:Nimbus和Supervisor被设计为了无状态的，所有的状态都位于Zookeeper和本地磁盘。只要用daemontools或monit做好进程的监控，失败就会自动重启。
- **Nimbus是否单点(SPOF)**:Nimbus失效后不会影响Supervisor监管Worker，但会影响Worker的重新分配。所以这里需要着重增加监管。

易于部署和运维

- 架构类似于MapReduce，便于运维人员维护

兼容多语言

Java、Scala、JRuby、Jython、Clojure..

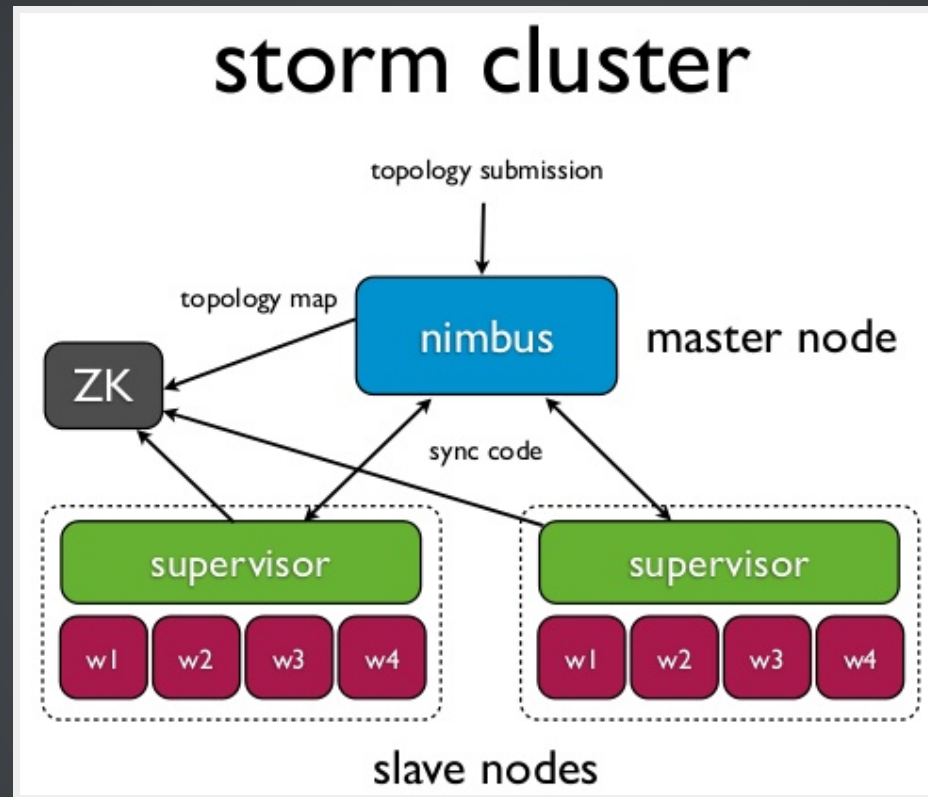
使用场景

-  流式处理

-  分布式RPC

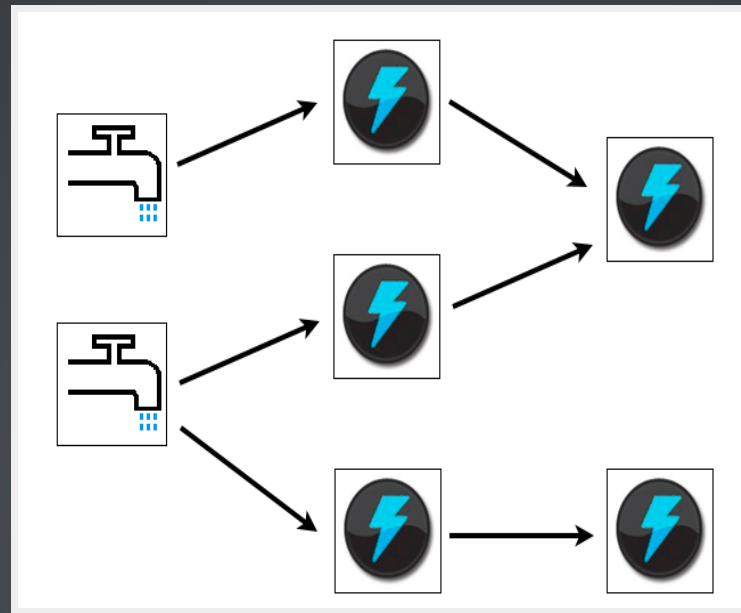
-  持续计算

系统构架



Hadoop	Storm
JobTracker	Nimbus
	Zookeeper
TaskTracker	Supervisor
Task	Worker

编程模型



Hadoop

Job

Map

Reduce

Writable

Task(Map/Reduce)

Storm

Topology

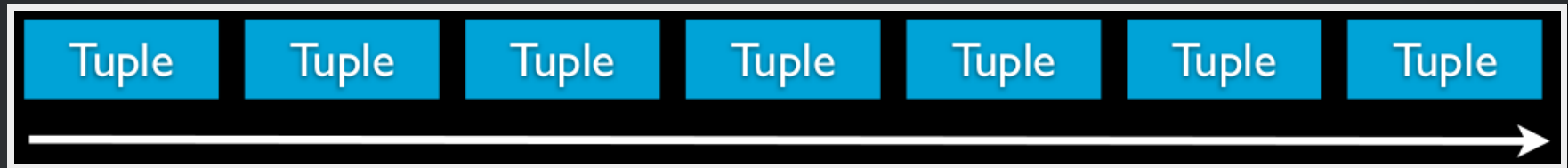
Spout

Bolt

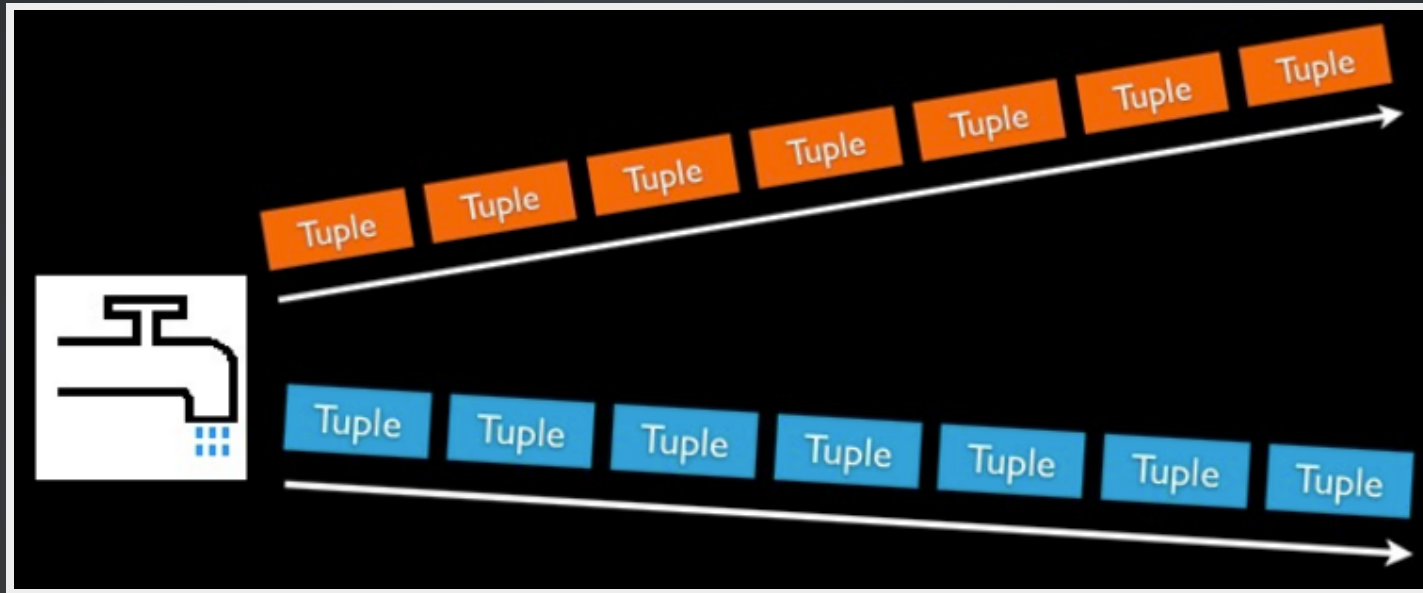
Tuple

Task(Spout/Bolt)

STREAMS

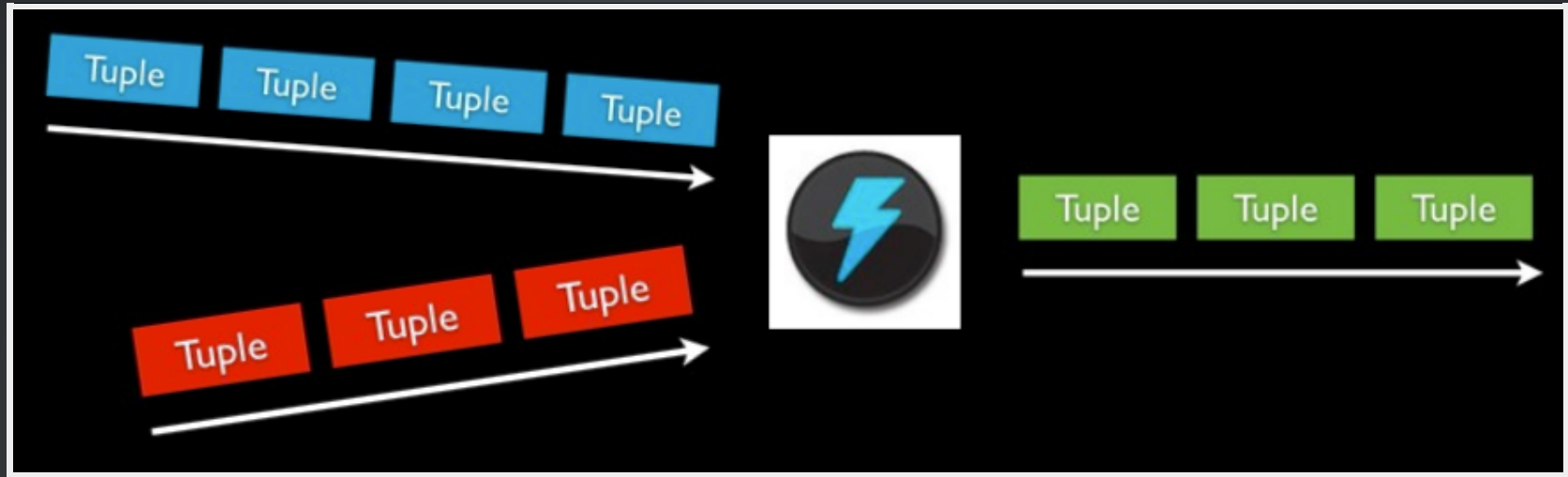


SPOUTS



```
public interface ISpout extends Serializable {  
    void open(Map conf, TopologyContext context,  
              SpoutOutputCollector collector);  
    void close();  
    void nextTuple();  
    void ack(Object msgId);  
    void fail(Object msgId);  
}
```

BOLTS

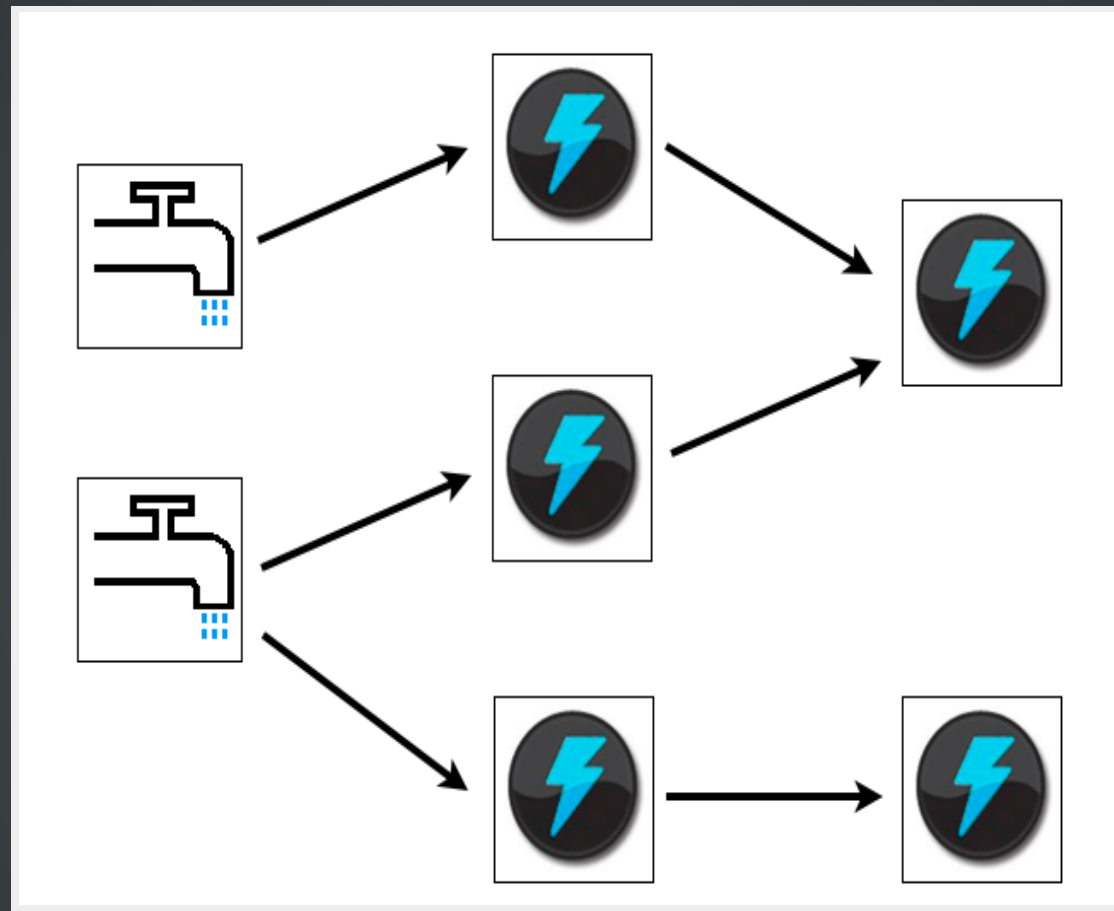


- 功能性操作
- 过滤操作
- 聚集操作
- Join操作
- 持久化操作(DB/Cache/HBase/Cassandra)

BOLTS

```
public interface IBolt extends Serializable {  
    void prepare(Map stormConf, TopologyContext context, OutputCollector  
    void execute(Tuple input);  
    void cleanup();  
}
```

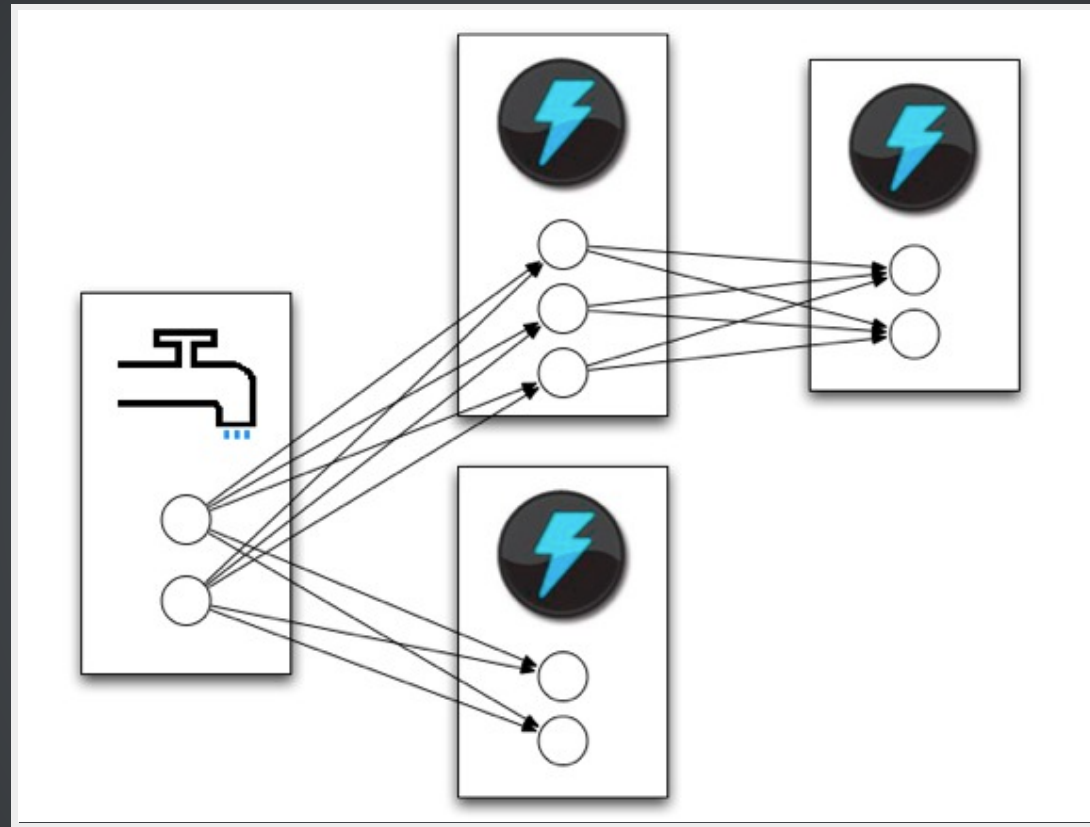
TOPOLOGIES



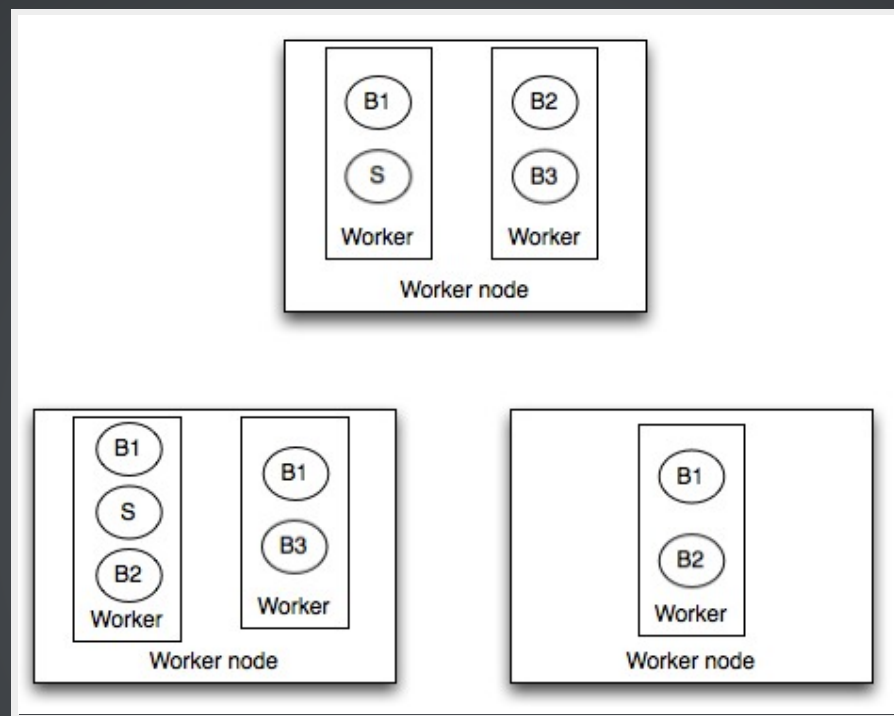
TOPOLOGIES

```
TopologyBuilder builder = new TopologyBuilder();
builder.setSpout("sentences", new KestrelSpout("kestrel.backtype.com",
                                                22133,
                                                "sentence_queue",
                                                new StringScheme()));
builder.setBolt("split", new SplitSentence(), 10)
    .shuffleGrouping("sentences");
builder.setBolt("count", new WordCount(), 20)
    .fieldsGrouping("split", new Fields("word"));
```

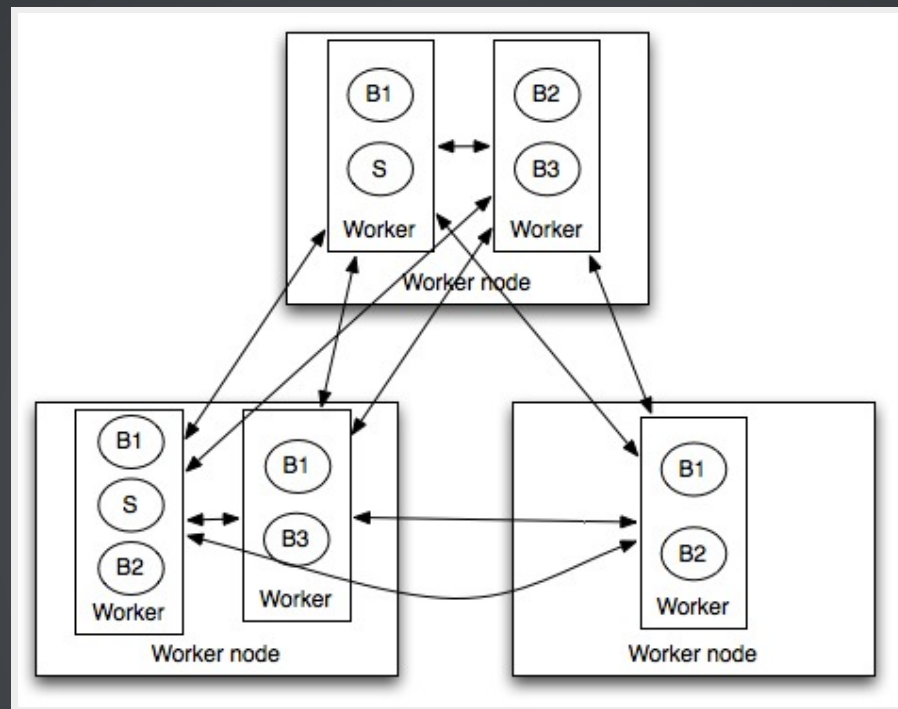
TASKS



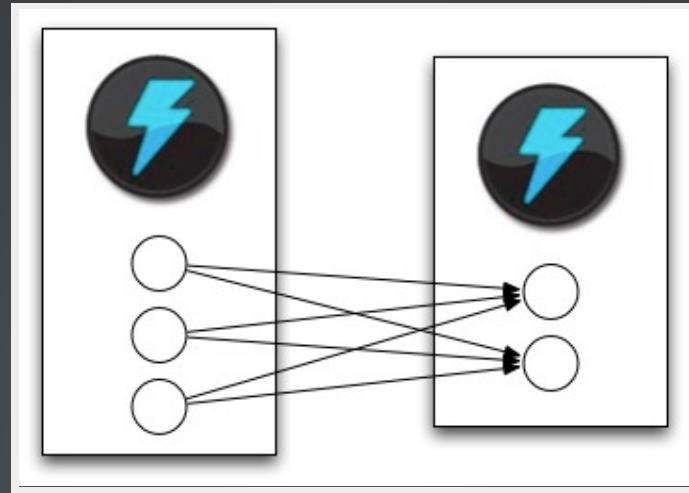
任务执行



任务执行

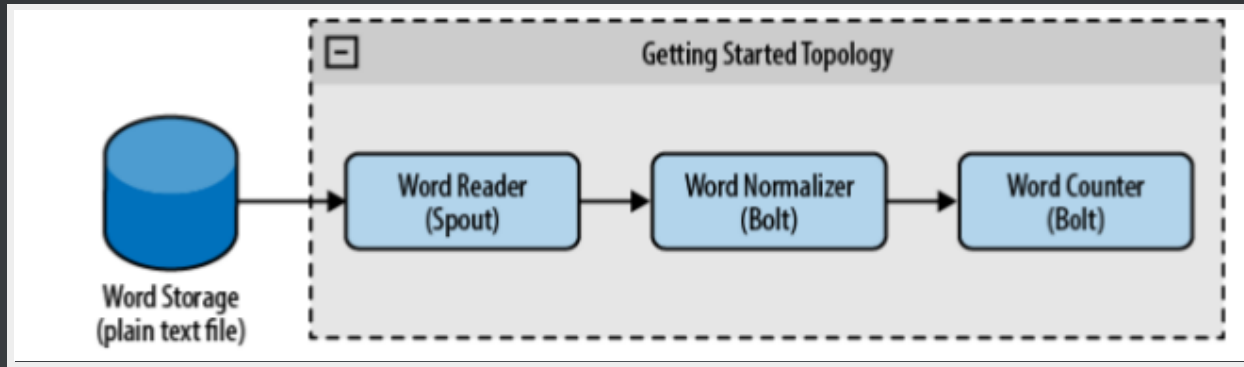


数据流分组



- Shuffle grouping:选择随机任务
- Fields Grouping:在tuple某字段上hash后求模
- All grouping:发送到所有的任务上
- Global grouping:取最小id的任务
- Direct:发送消息时指定任务
- Custorm:用户自定义,实现CustormStreamGrouping

WORD COUNT



WORDREADER

```
public class WordReader implements IRichSpout {
    private SpoutOutputCollector collector;
    private FileReader fileReader;
    private boolean completed = false;
    private TopologyContext context;
    public boolean isDistributed() {return false;}
    public void ack(Object msgId) {
        System.out.println("OK:"+msgId);
    }
    public void close() {}
    public void fail(Object msgId) {
        System.out.println("FAIL:"+msgId);
    }
    /**
     * The only thing that the methods will do It is emit each
     * file line
     */
}
```

WORDNORMALIZER

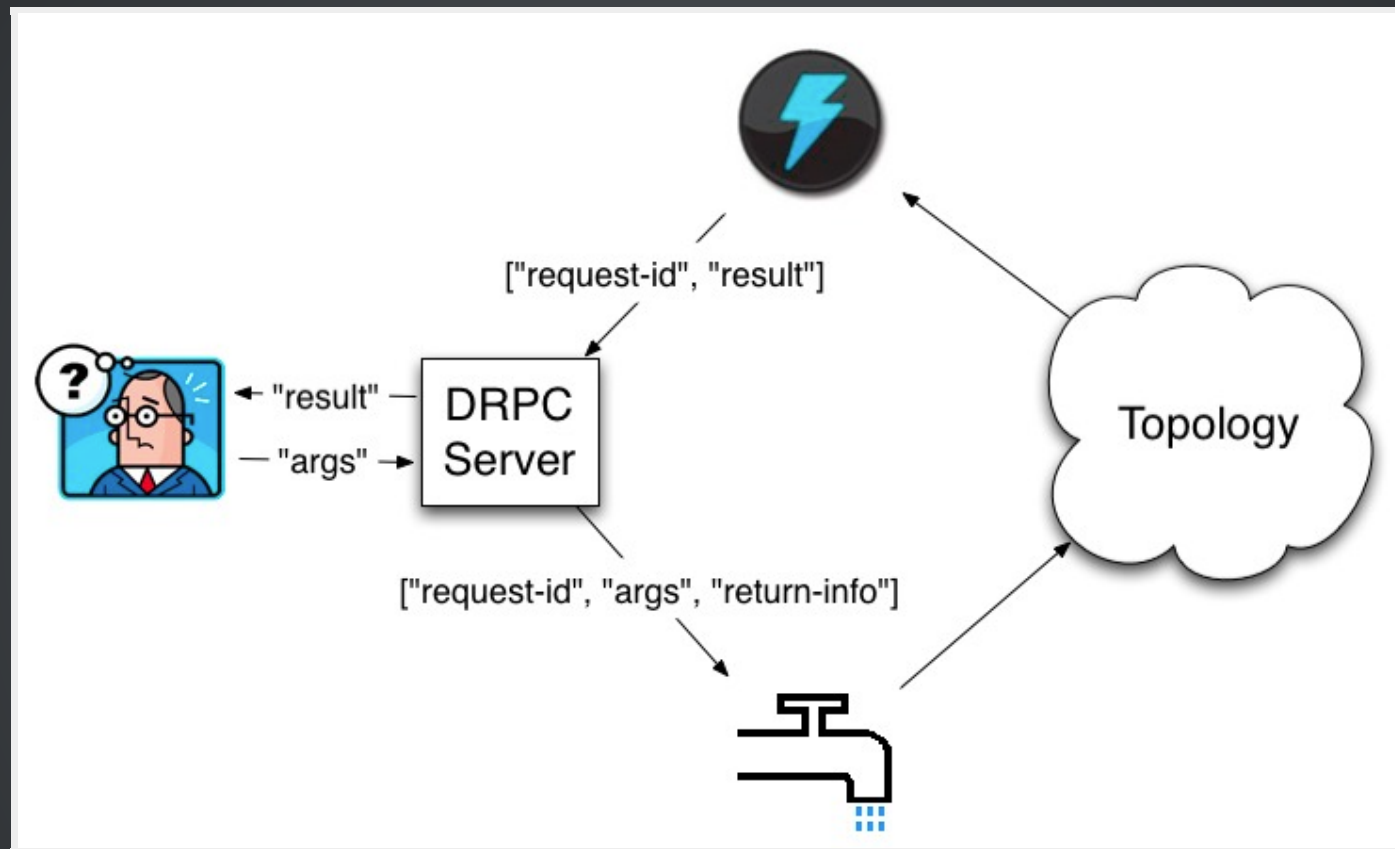
```
public class WordNormalizer implements IRichBolt {
    private OutputCollector collector;
    public void cleanup() {}
    /**
     * The bolt will receive the line from the
     * words file and process it to Normalize this line
     *
     * The normalize will be put the words in lower case
     * and split the line to get all words in this
     */
    public void execute(Tuple input) {
        String sentence = input.getString(0);
        String[] words = sentence.split(" ");
        for(String word : words){
            word = word.trim();
            if(!word.isEmpty()){
                word = word.toLowerCase();
            }
        }
    }
}
```

WORDCOUNTER

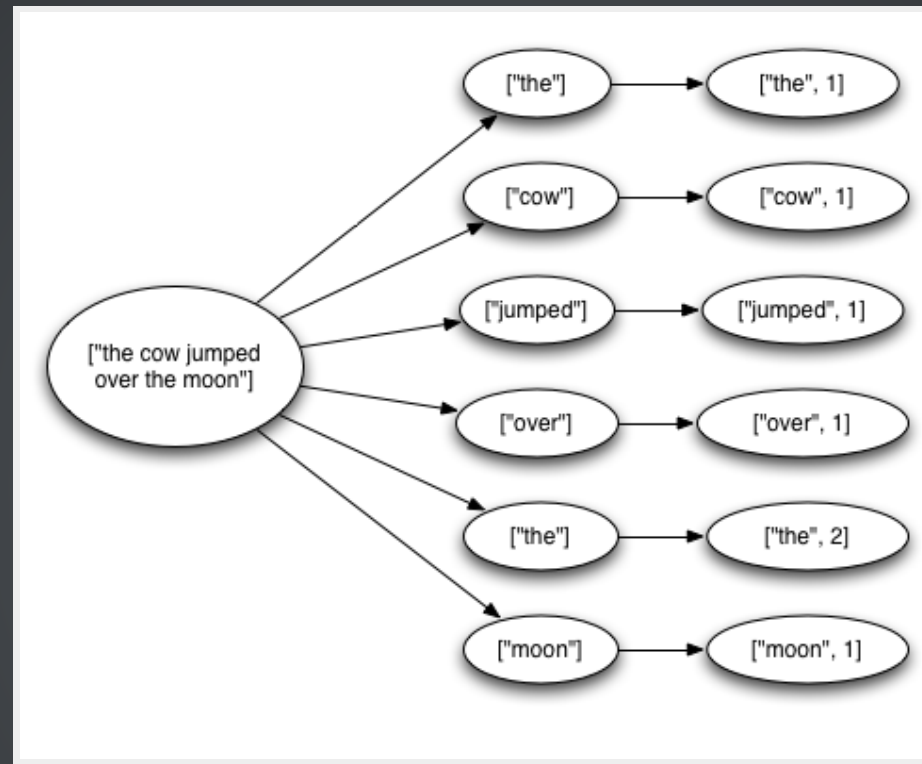
```
public class WordCounter implements IRichBolt {
    Integer id;
    String name;
    Map<string, integer=""> counters;
    private OutputCollector collector;

    /**
     * At the end of the spout (when the cluster is shutdown
     * We will show the word counters
     */
    @Override
    public void cleanup() {
        System.out.println("-- Word Counter ["+name+"-"+id+"] --");
        for(Map.Entry<string, integer=""> entry : counters.entrySet()){
            System.out.println(entry.getKey()+" : "+entry.getValue());
        }
    }
}
```

DRPC



记录级容错



记录级容错

- Tuple超过topology.message.timeout.secs未被处理被视为失效
- 首先Storm从Spout中请求nextTuple获取tuple
- Spout通过SpoutOutputCollector提供的open方法输出tuple到输出流,并提供一个message id标识此tuple
- 接着tuple被Bolt消费,Storm跟踪它创建的消息树。如果Storm检测到tuple被fully process , Storm调用ack传入该message id。如果tuple超时将在Spout中调用fail方法。

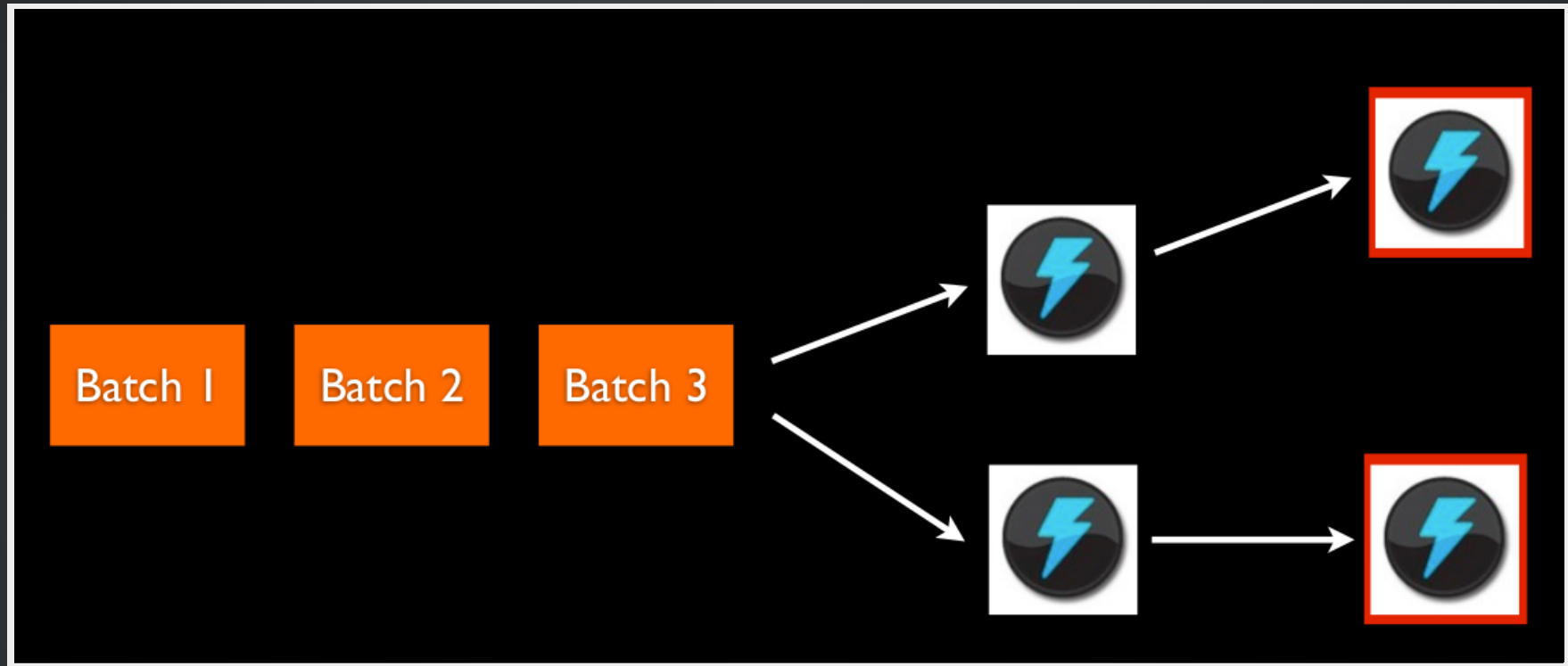
事务拓扑

- Transactional topology
- Trident

TRANSACTIONAL TOPOLOGY

- 0.7.0引入
- 需要支持指定批次replay的queue
- 实现了基于kafka的transactional spout

TRANSACTIONAL TOPOLOGY



TRIDENT

the cow jumped over the moon
the man went to the store and bought some candy
four score and seven years ago
how many apples can you eat
the cow jumped over the moon
the man went to the store and bought some candy
four score and seven years ago
how many apples can you eat
the cow jumped over the moon
the man went to the store and bought some candy



the cow jumped over the moon
the man went to the store and bought some candy
four score and seven years ago
Batch 1

how many apples can you eat
the cow jumped over the moon
the man went to the store and bought some candy
four score and seven years ago
how many apples can you eat
Batch 2

the cow jumped over the moon
the man went to the store and bought some candy
Batch 3

- 0.8.0引入

TRIDENT SPOUT

```
FixedBatchSpout spout = new FixedBatchSpout(new Fields("sentence"), 3,  
    new Values("the cow jumped over the moon"),  
    new Values("the man went to the store and bought some candy"),  
    new Values("four score and seven years ago"),  
    new Values("how many apples can you eat"),  
    spout.setCycle(true);
```

TRIDENT TOPOLOGY

```
TridentTopology topology = new TridentTopology();
TridentState wordCounts =
    topology.newStream("spout1", spout)
        .each(new Fields("sentence"), new Split(), new Fields("word"))
        .groupBy(new Fields("word"))
        .persistentAggregate(new MemoryMapState.Factory(), new Count(), ne
        .parallelismHint(6);
```

TRIDENT BASEFUNCTION

```
public class Split extends BaseFunction {  
    public void execute(TridentTuple tuple, TridentCollector collector) {  
        String sentence = tuple.getString(0);  
        for(String word: sentence.split(" ")) {  
            collector.emit(new Values(word));  
        }  
    }  
}
```

TRIDENT DRPC

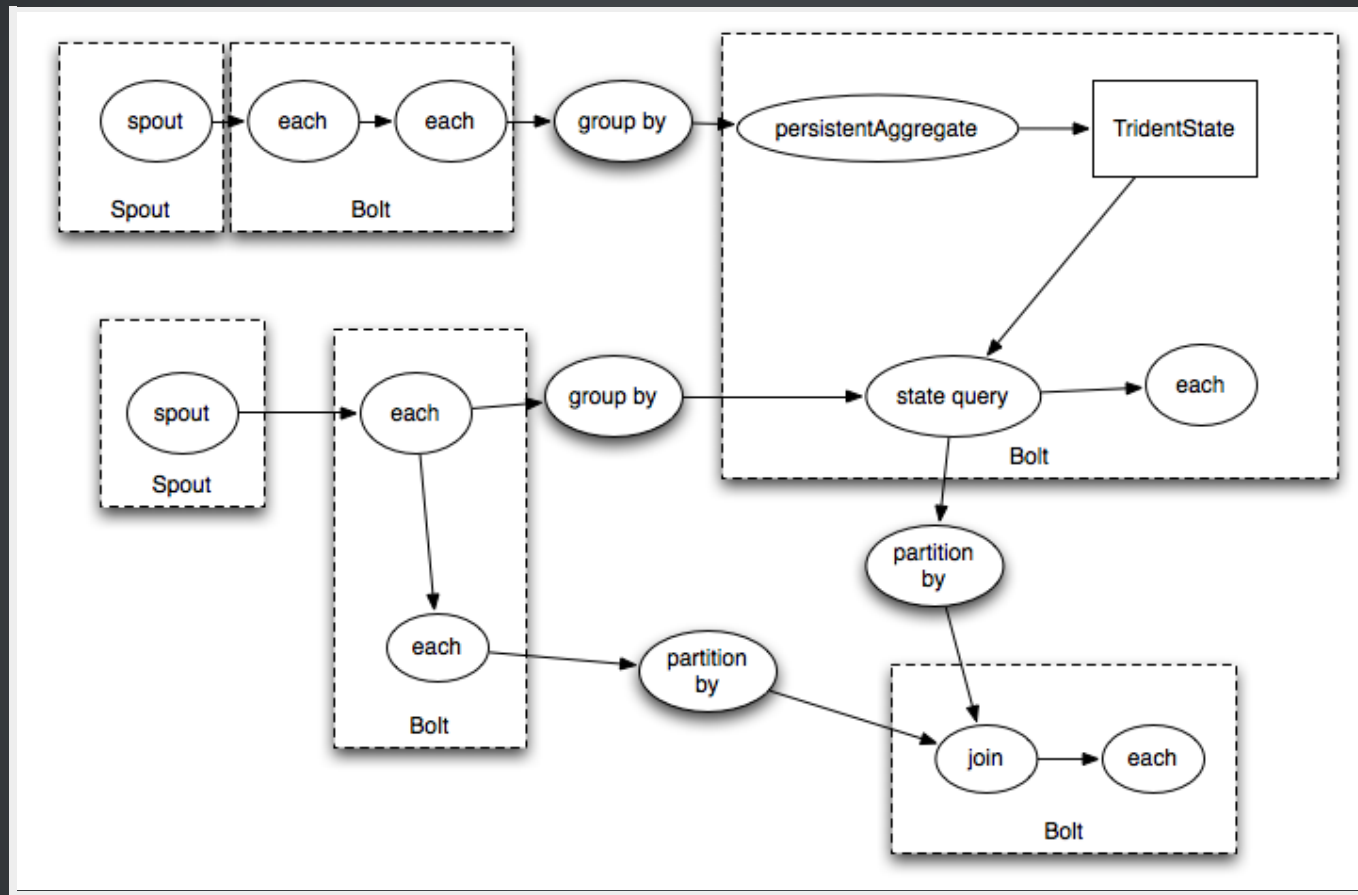
```
DRPCClient client = new DRPCClient("drpc.server.location", 3772);  
System.out.println(client.execute("words", "cat dog the man");  
// prints the JSON-encoded result, e.g.: "[[5078]]"
```

```
topology.newDRPCStream("words")  
    .each(new Fields("args"), new Split(), new Fields("word"))  
    .groupBy(new Fields("word"))  
    .stateQuery(wordCounts, new Fields("word"), new MapGet(), new Fields("count"))  
    .each(new Fields("count"), new FilterNull())  
    .aggregate(new Fields("count"), new Sum(), new Fields("sum"));
```


状态存储

- Trident解决的两个问题
 - 为每一批次提供了一个唯一的transaction id，重试时使用一致的transaction id
 - 状态更新依据批次的顺序
- 存储方式
 - Memory
 - Memcached
 - Cassandra
 - ...

TRIDENT TOPOLOGY



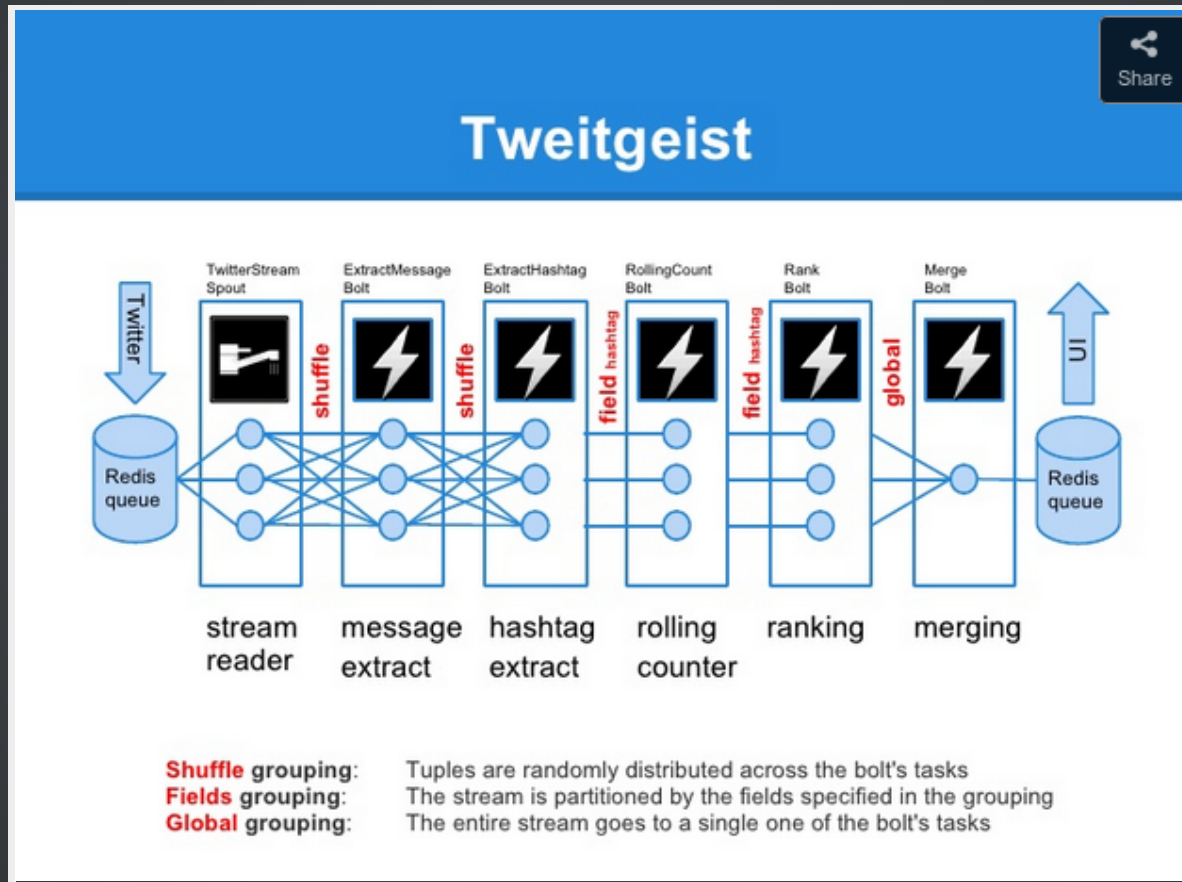
STORM的问题

- 向topology增加业务时需要删除原有topology部署新的topology(0.9增加了swap)
- topology间不能传递数据（只能通过Queue队列间接传递）

业界应用

TWITTER场景(TWEITGEIST)

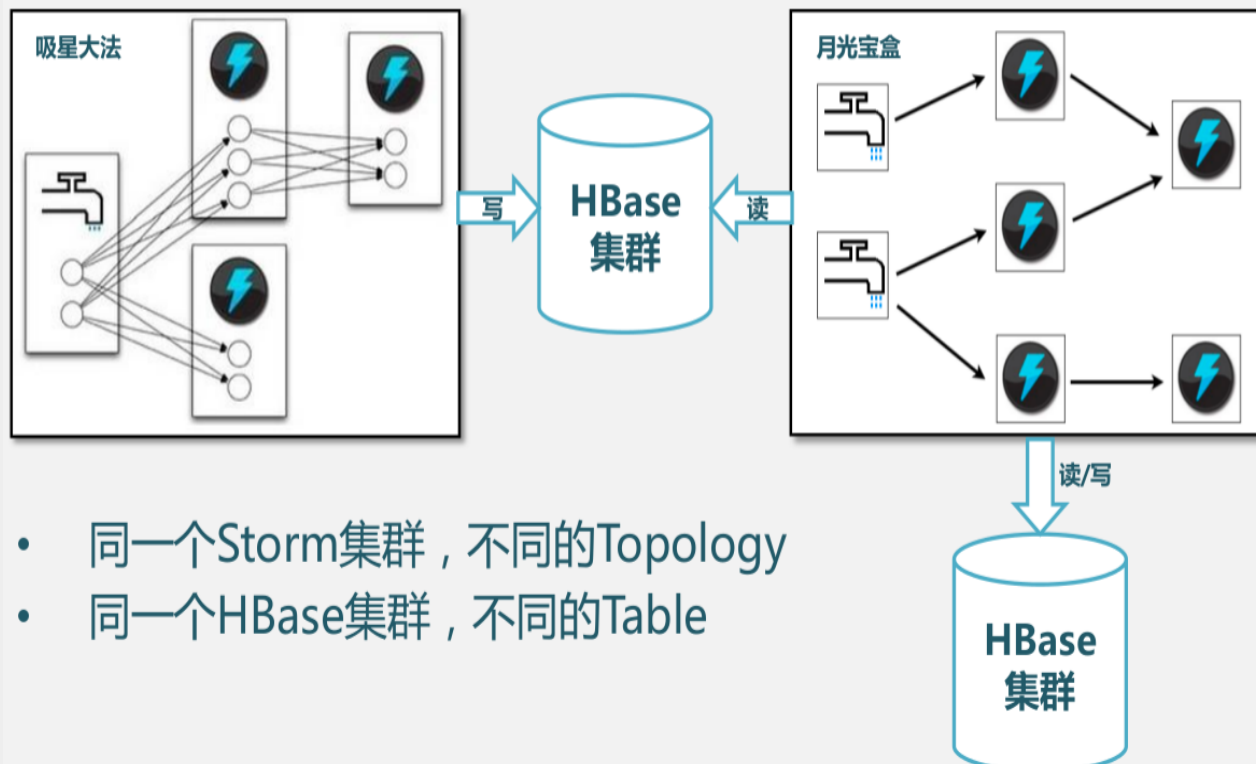
LIVE TOP 10 TRENDING HASHTAGS ON TWITTER



一淘数据部(月光宝盒)

月光宝盒项目简介

一淘数据部
etao.com
数据让生活更简单!

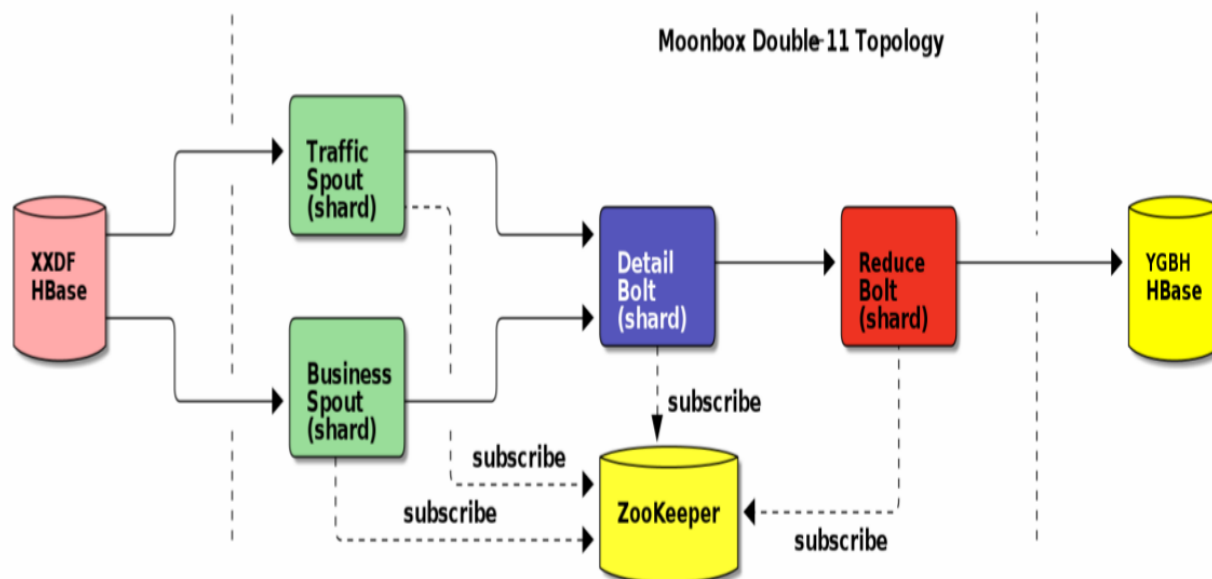


一淘数据部(月光宝盒)

月光宝盒项目简介

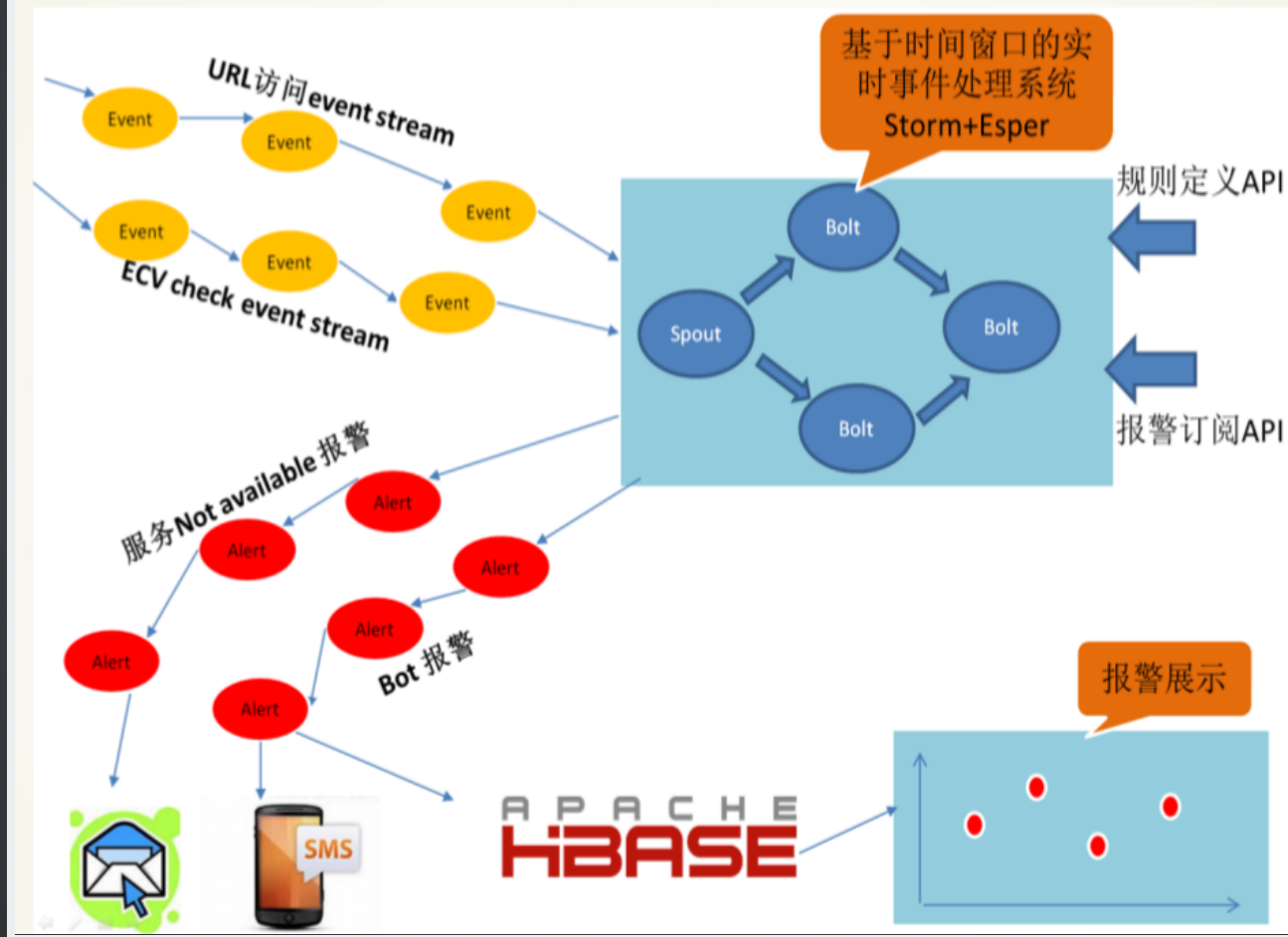
一淘数据部
etao.com
数据让生活更简单!

- 双11拓扑任务



携程(ALERTING)

Alerting - Framework



运维监控

监控工具

- Storm UI
- JMX/VisualVM
- Yammer Metrics
- Ganglia/Graphite
- Log4j
- Nagios(监控硬件、系统、log中的ERROR/WARN)
- OpenTSDB
- storm-monitor

STORM UI

JMX/VISUALVM

YAMMER METRICS

[HTTP://METRICS.CODAHALE.COM/](http://metrics.codahale.com/)

```
<dependencies>
  <dependency>
    <groupId>com.codahale.metrics</groupId>
    <artifactId>metrics-core</artifactId>
    <version>${metrics.version}</version>
  </dependency>
</dependencies>
```

YAMMER METRICS

```
//application level static
final MetricRegistry metrics = new MetricRegistry();
//Gauges 标尺
public class QueueManager {
    private final Queue queue;

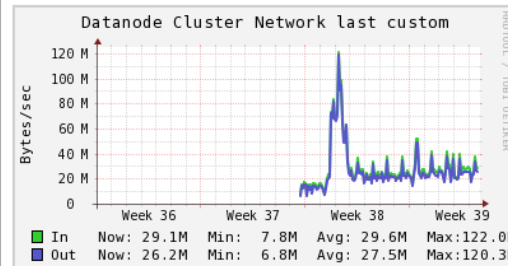
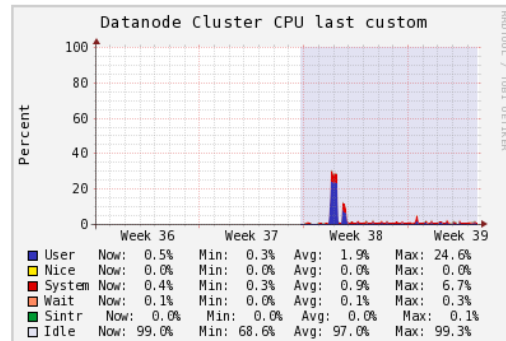
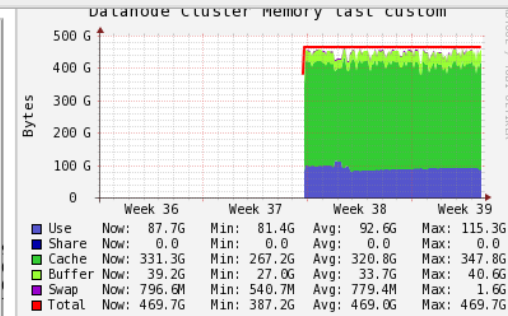
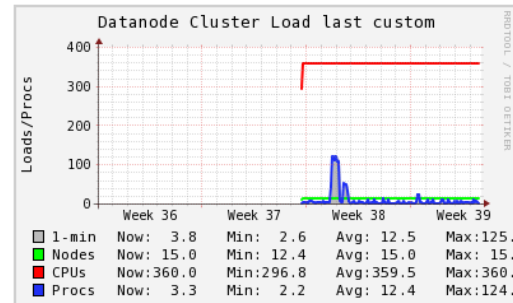
    public QueueManager(MetricRegistry metrics, String name) {
        this.queue = new Queue();
        metrics.register(MetricRegistry.name(QueueManager.class, name, "s
            new Gauge<integer>() {
                @Override
                public Integer getValue() {
                    return queue.size();
                }
            });
    }
}
```

GANGLIA

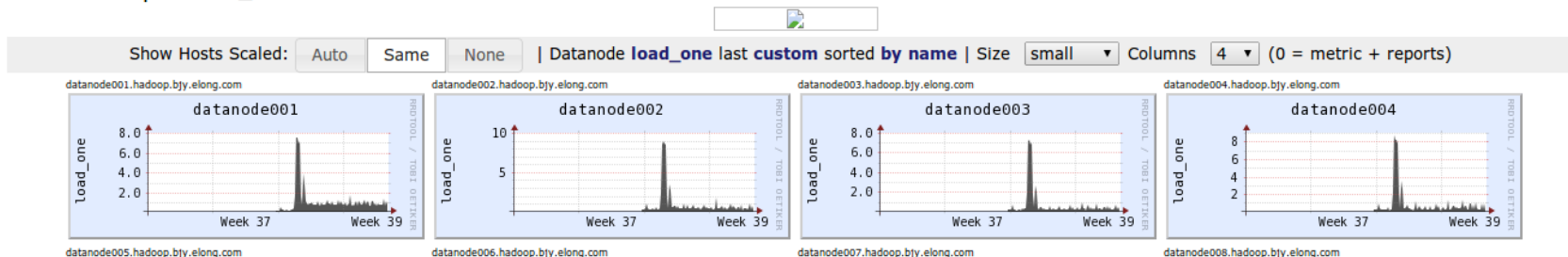
Hosts up: **16**
Hosts down: **0**

Current Load Avg (15, 5, 1m):
1%, 1%, 1%
Avg Utilization (last custom):
1%

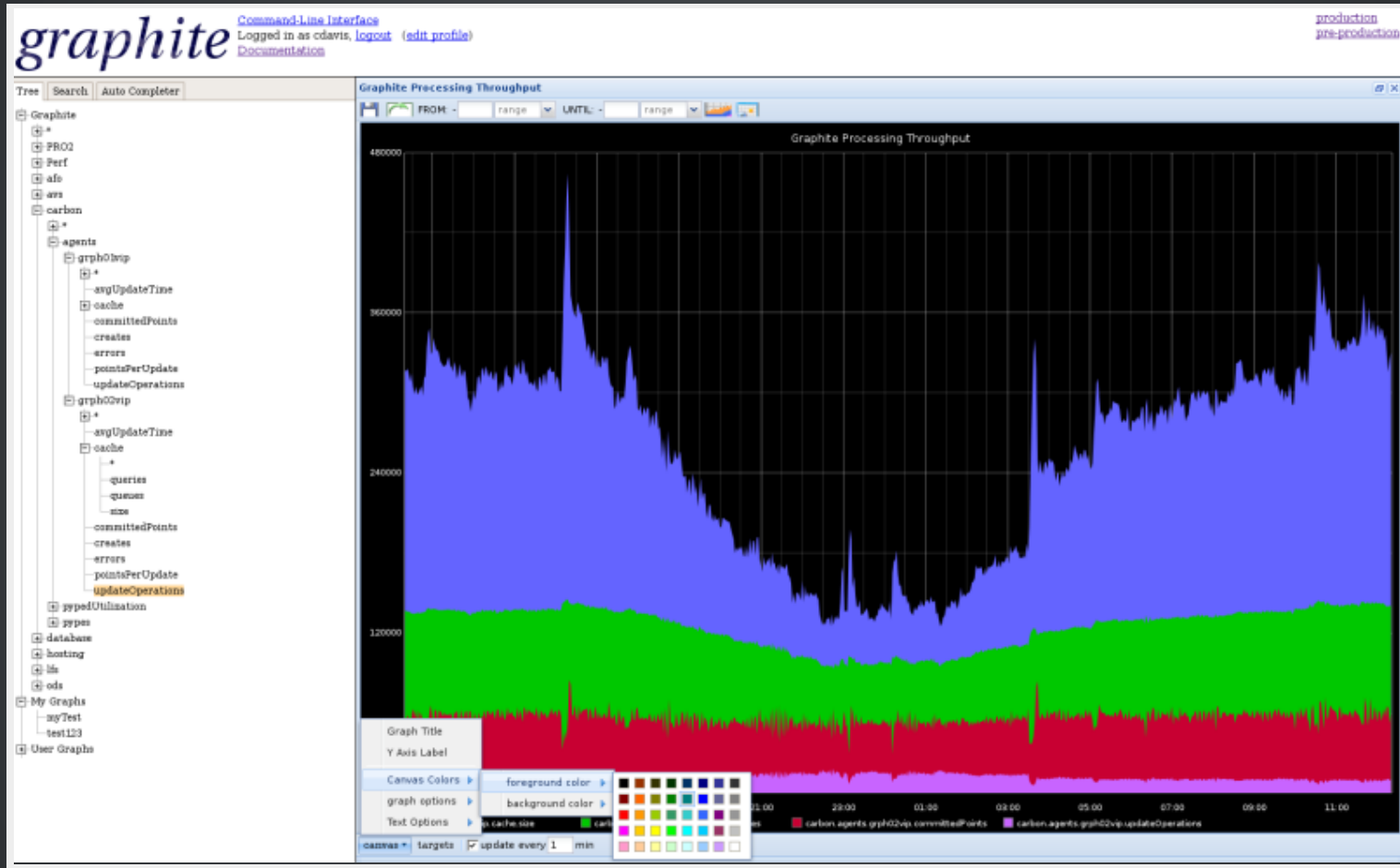
Utilization heatmap



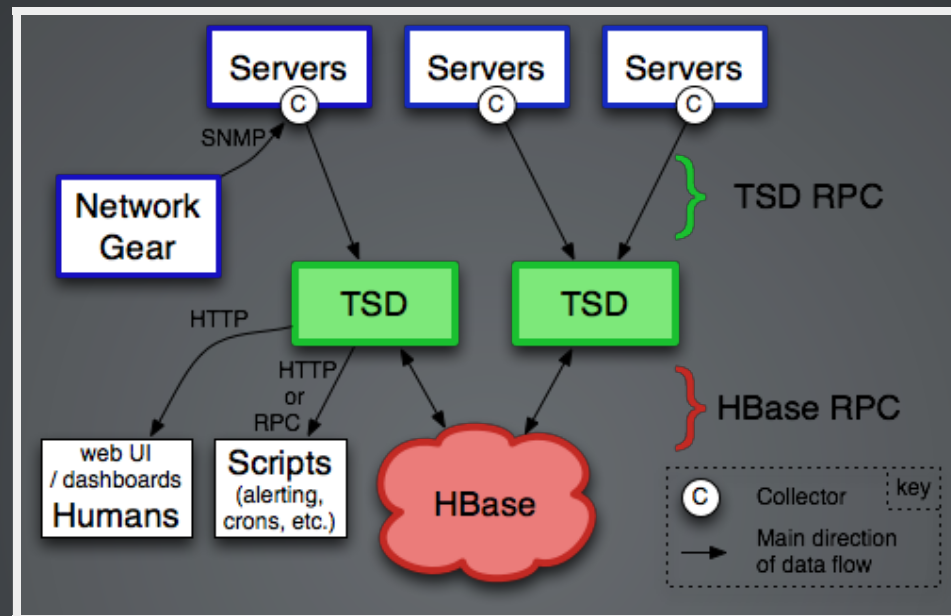
Stacked Graph - load_one

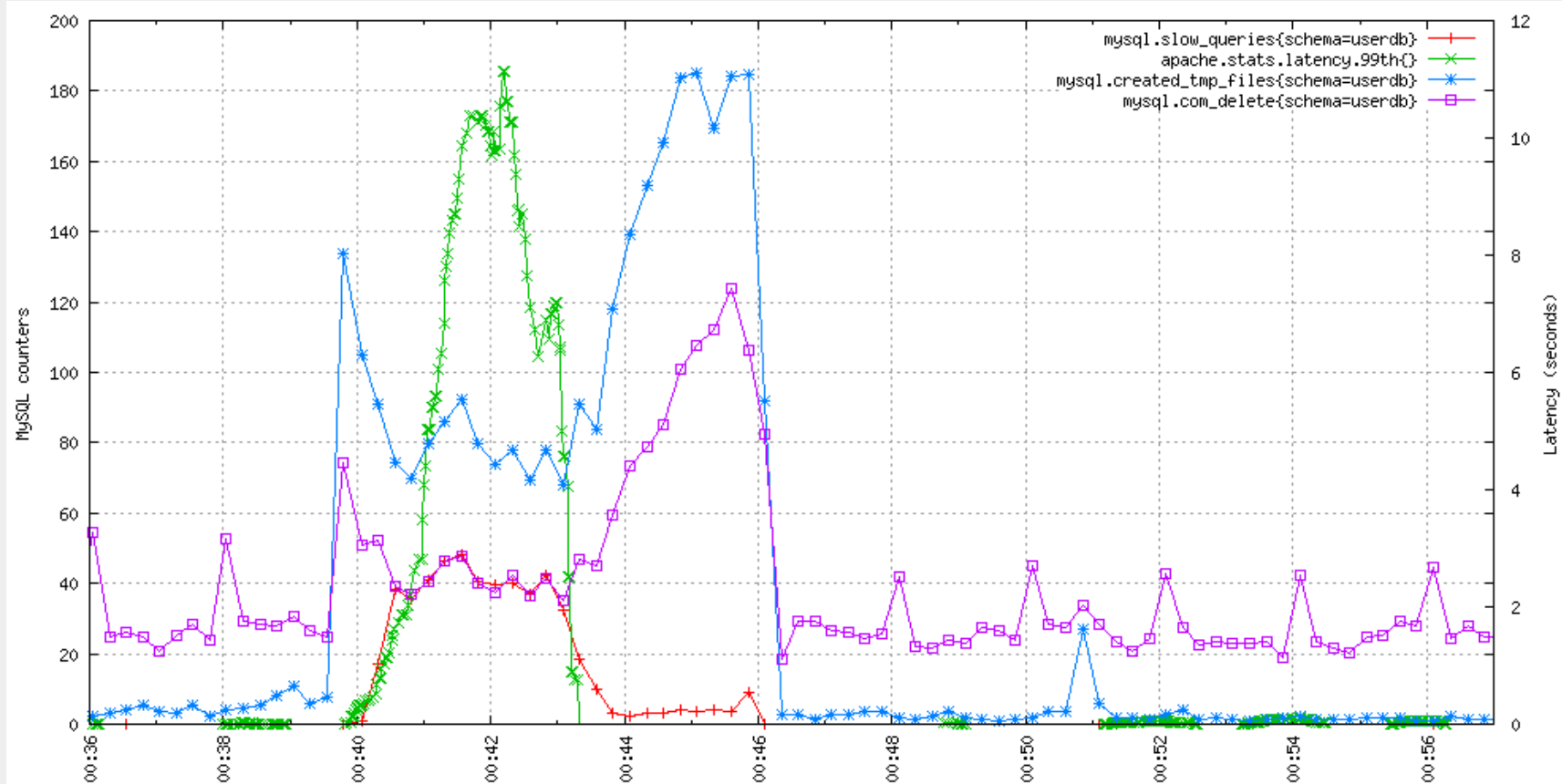


GRAPHITE



OPENTSDB





STORM-MONITOR

[HTTPS://GITHUB.COM/KILLME2008/STORM-MONITOR](https://github.com/KILLME2008/STORM-MONITOR)参考用于监控NIMBUS

- 监控supervisor数目是否正确，当supervisor挂掉的时候会发送警告
- 监控nimbus是否正常运行，monitor会尝试连接nimbus，如果连接失败就认为nimbus挂掉。
- 监控topology是否正常运行，包括它是否正常部署，是否有运行中的任务

其它可选监控工具

- Ooyala metrics_storm
- Storm 0.9 Metrics
 - 按固定的时间窗口收集任意的用户自定义metric
 - metric最终被Storm聚合
 - 可自定义MetricsConsumer

服务性能监控

- 拓扑图
 - 吞吐量
 - 延迟
- Bolt和spout
 - 吞吐量
 - 延迟
 - 队列时间
 - 执行时间
- 数据来源和数据存储
 - 读/写条数
 - 读/写速度
- Storm 0.8.2 new metrics
 - process latency(tuple直到被ack的时间)
 - execute latency(tuple被调度执行的时间)
 - capacity(上十分钟bolt用于执行tuple的时间百分比)

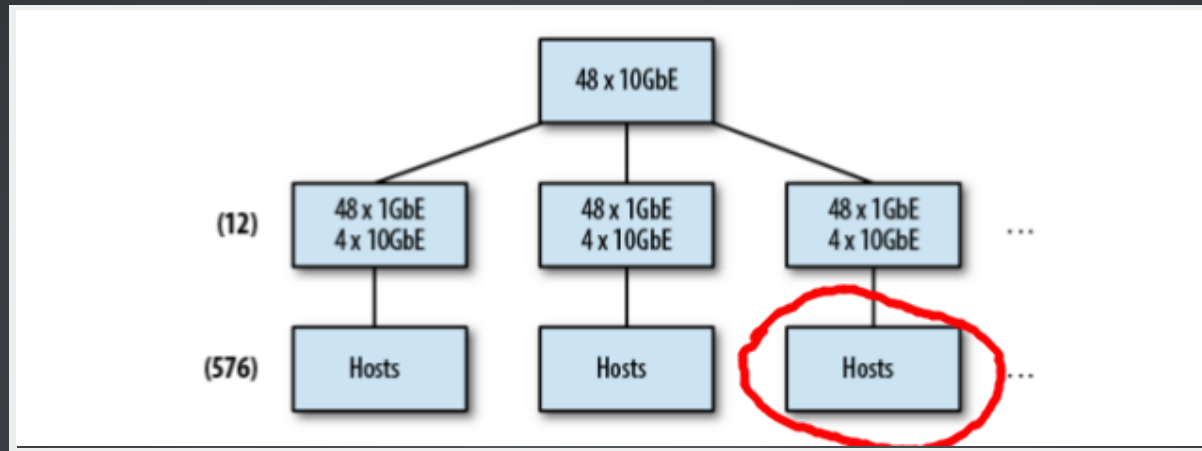
CHECKLIST应用级监控

- Kafka读取速度
- HBase写入速度
- 规则处理引擎某topic单条数据处理时间
- 规则处理引擎处理速度
- 规则处理引擎出错次数
- 处理数据延迟(根据数据业务日期)
- QueryTools查询响应时间
- 数据接口监控

推荐监控方案

- 服务或组件监控
 - Nagios(硬件、服务进程、log)参考storm_monitor
 - Storm UI(Thrift)
 - Ooyala metrics_storm
 - Ganglia/Graphite
- 应用监控
 - Yammer Metrics

网络规划



THE END